



Green Destiny: A 240-Node Compute Cluster in One Cubic Meter

Wu-chun (Wu) Feng

Research & Development in Advanced Network Technology (RADIANT)
Computer & Computational Sciences Division
Los Alamos National Laboratory



Outline

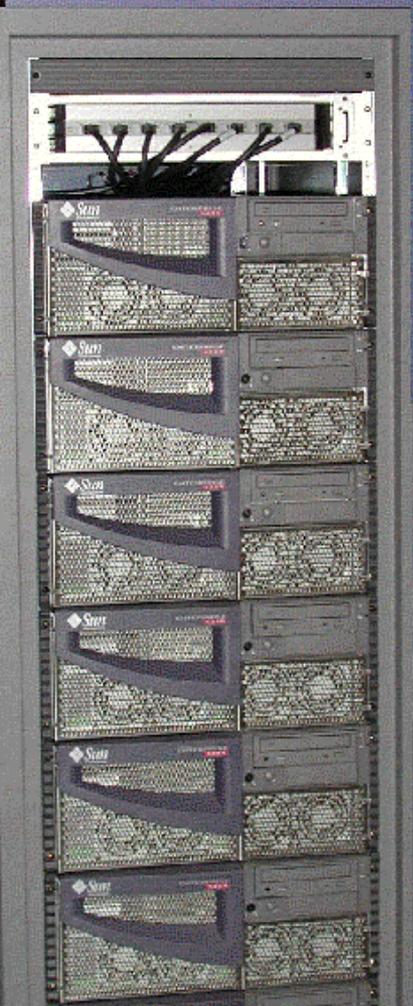
- Where is Supercomputing?
 - ◆ Architectures from the Top 500.
- Evaluating Supercomputers
 - ◆ Metrics: Performance & Price/Performance
- An Alternative Flavor of Supercomputing
 - ◆ Supercomputing in Small Spaces → Bladed Beowulf
- Architecture of a Bladed Beowulf
- Performance Metrics
- Benchmark Results
- Discussion & Status
- Conclusion
- Acknowledgements & Media Coverage



Flavors of Supercomputing

(Picture Source: Thomas Sterling, Caltech & NASA JPL)

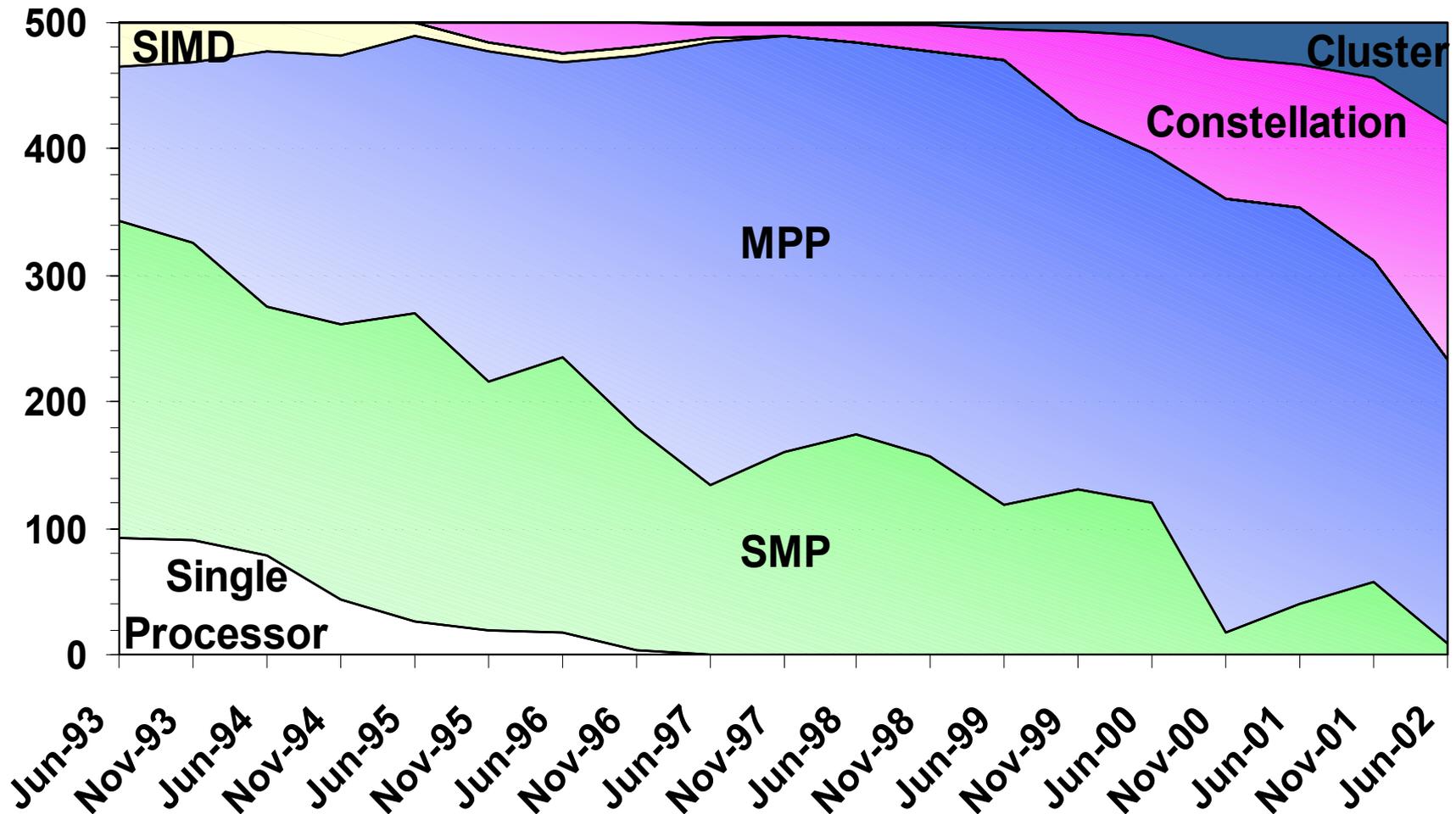
Sun Microsystems, Inc.
Myrinet Technical Compute Farm



COMPAQ AlphaServer SC



Architectures from the Top 500 Supercomputer List





Metrics for Evaluating Supercomputers

- *Performance*
 - ◆ Metric: Floating-Operations Per Second (FLOPS)
 - ◆ Example: Japanese Earth Simulator
- *Price/Performance → Cost Efficiency*
 - ◆ Metric: Cost / FLOPS
 - ◆ Examples: SuperMike, GRAPE-5, Avalon.



Performance (At Any Cost)

- Japanese Earth Simulator (\$400M)

	Performance	Price/Perf
Peak	40.00 Tflop	\$10.00/Mflop
Linpack	35.86 Tflop	\$11.15/Mflop
n-Body	29.50 Tflop	\$13.56/Mflop
Climate	26.58 Tflop	\$15.05/Mflop
Turbulence	16.40 Tflop	\$24.39/Mflop
Fusion	14.90 Tflop	\$26.85/Mflop



Price/Performance

■ Cost Efficiency

- ◆ LSU's SuperMike
(2002: \$2.8M)

	Performance	Price/Perf
Linpak	2210 Gflops	\$1.27/Mflop

- ◆ U. Tokyo's GRAPE-5
(1999: \$40.9K)

	Performance	Price/Perf
N-body	5.92 Gflops	\$6.91/Mflop

- ◆ LANL's Avalon
(1998: \$152K)

	Performance	Price/Perf
Peak	149.40 Gflops	\$1.02/Mflop
Linpak	19.33 Gflops	\$7.86/Mflop



The Need for *New* Supercomputing Metrics

- Analogy: Buying a car. Which metric to use?
 - ◆ Raw performance, price/performance, fuel efficiency, reliability, size, etc.
- Issues with today's supercomputing metrics
 - ◆ Focus: Performance & price/performance
 - ☞ Important metrics, but ...



Flavors of Supercomputing

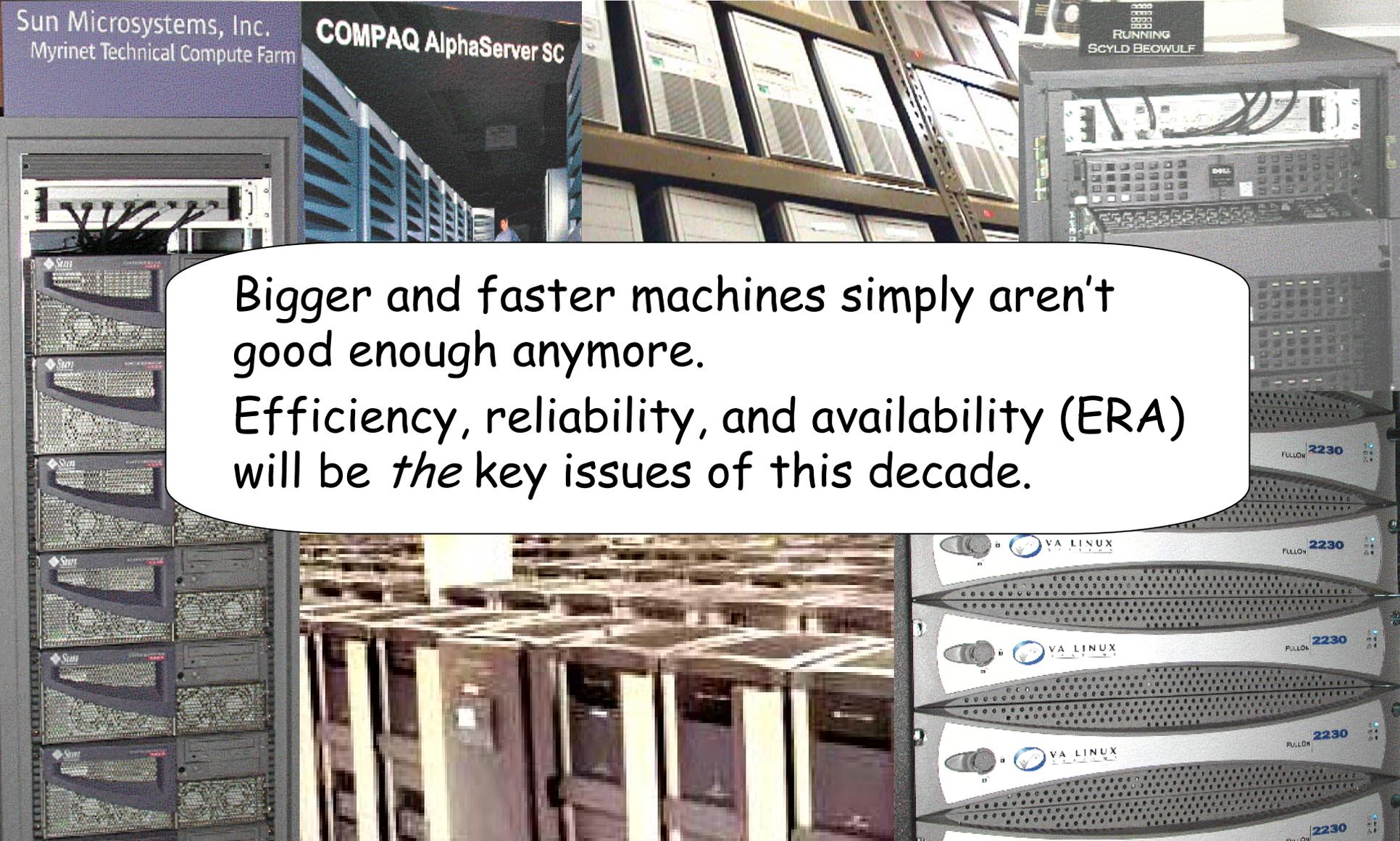
(Picture Source: Thomas Sterling, Caltech & NASA JPL)

Sun Microsystems, Inc.
Myrinet Technical Compute Farm

COMPAQ AlphaServer SC

RUNNING
SCYLD BEOWULF

Bigger and faster machines simply aren't good enough anymore. Efficiency, reliability, and availability (ERA) will be *the* key issues of this decade.





Why ERA Metrics?

■ Observations

- ◆ Strong hints of the tradeoffs that come with "performance" and "price/performance" metrics ...
 - ☞ Lower efficiency, reliability, and availability.
 - ☞ Higher operational costs, e.g., admin, maintenance, etc.
- ◆ Institutional consumers that use clusters as a tool ...
 - ☞ Pharmaceutical, financial, actuarial, retail, aerospace, data centers for web-server farms.
 - ☞ A couple of informational data points:
 - Peter Bradley, Pratt & Whitney: IEEE Cluster 2002.
 - Reliability, transparency, and ease of use.
 - Eric Schmidt, Google: IEEE Hot Chips & NY Times, 2002.
 - Low power, NOT speed.
 - DRAM density, NOT speed.



An Alternative Flavor of Supercomputing

- Supercomputing in Small Spaces (<http://sss.lanl.gov>)
 - ◆ First instantiation: *Bladed Beowulf*
 - ☞ MetaBlade (24), MetaBlade2 (24), and Green Destiny (240).
- Goal
 - ◆ Improve *efficiency, reliability, and availability* (ERA) in large-scale computing systems.
 - ☞ Sacrifice a little bit of raw performance.
 - ☞ Improve overall system throughput as the system will “always” be available, i.e., effectively no downtime, no hardware failures, etc.
 - ◆ Reduce the *total cost of ownership* (TCO).
- Analogy
 - ◆ Ferrari 550: Wins raw performance but reliability is poor so it spends its time in the shop. Throughput low.
 - ◆ Toyota Camry: Loses raw performance but high reliability results in high throughput (i.e., miles driven).



Architecture of a Bladed Beowulf

A Fundamentally Different Approach to
High-Performance Computing



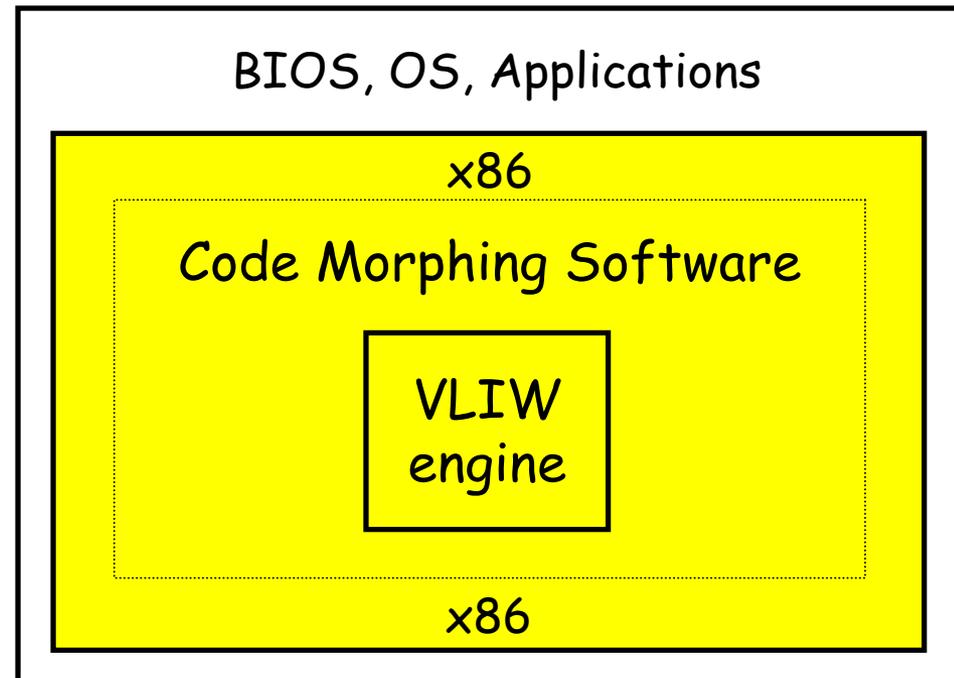
Transmeta TM5600 CPU: VLIW + CMS

■ VLIW Engine

- ◆ Up to four-way issue
 - ☞ In-order execution only.
 - ☞ 20% reduction on transistor count w.r.t superscalar arch.
- ◆ Two integer units
- ◆ Floating-point unit
- ◆ Memory unit
- ◆ Branch unit

■ VLIW Transistor Count ("Anti-Moore's Law")

- ◆ $\sim \frac{1}{4}$ of Intel PIII $\rightarrow \sim 6x-7x$ less power dissipation
- ◆ Less power \rightarrow lower "on-die" temp. \rightarrow better reliability & availability





Transmeta TM5x00 Comparison

Intel P4	MEM	MEM	2xALU	2xALU	FPU	SSE	SSE	Br
Transmeta TM5x00	MEM		2xALU		FPU			Br

- Current-generation Transmeta TM5800 performs comparably to an Intel PIII over iterative scientific codes on a clock-for-clock-cycle basis.
- Next-generation Transmeta CPU rectifies the above mismatch in functional units.

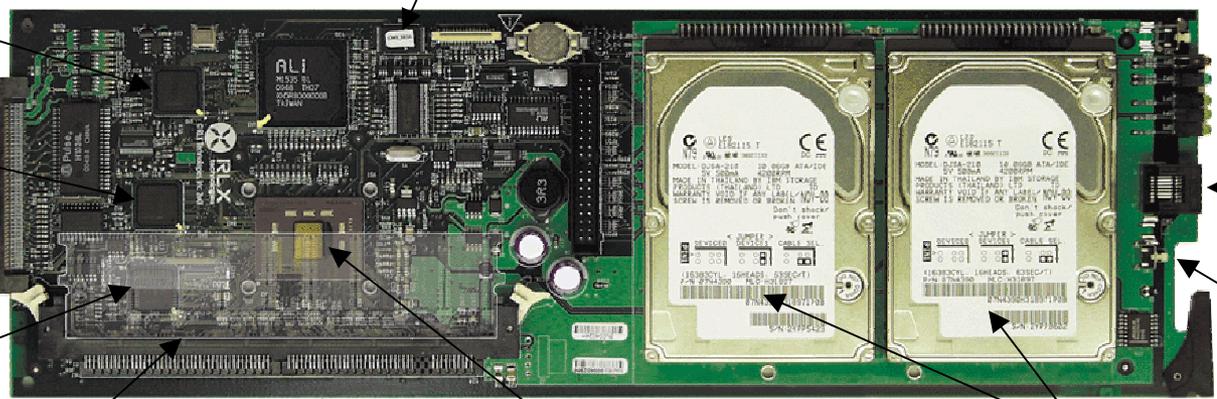


Transmeta TM5x00 CMS

- Code Morphing Software (CMS)
 - ◆ Provides compatibility by dynamically “morphing” x86 instructions into simple VLIW instructions.
 - ◆ Learns and improves with time, i.e., iterative execution.
- Modules for CMS
 - ◆ Interpreter
 - ☞ Interprets x86 instructions (*a la* Java).
 - ☞ Filters infrequently executed code from being optimized.
 - ☞ Collects run-time statistical information.
 - ◆ Translator
 - ☞ Re-compiles x86 instructions into optimized VLIW instructions (*a la* JIT compiler).



RLX ServerBlade™ 633 (circa 2000)



Code Morphing Software (CMS), 1 MB

Public NIC
33 MHz PCI

Status LEDs

Private NIC
33 MHz PCI

Serial RJ-45 debug port

Management NIC
33 MHz PCI

Reset Switch

128MB, 256MB, 512MB
DIMM SDRAM
PC-133

512KB
Flash ROM

Transmeta™
TM5600 633 MHz

ATA 66
0 or 1 or 2 - 2.5" HDD
10 or 30 GB each



128KB L1 cache, 512KB L2 cache
LongRun, Northbridge, x86 compatible

RLX ServerBlade™ 667 @ \$960.
933 @ TBD.
1066 @ alpha.



RLX System™ 324

3U chassis that houses 24 blades



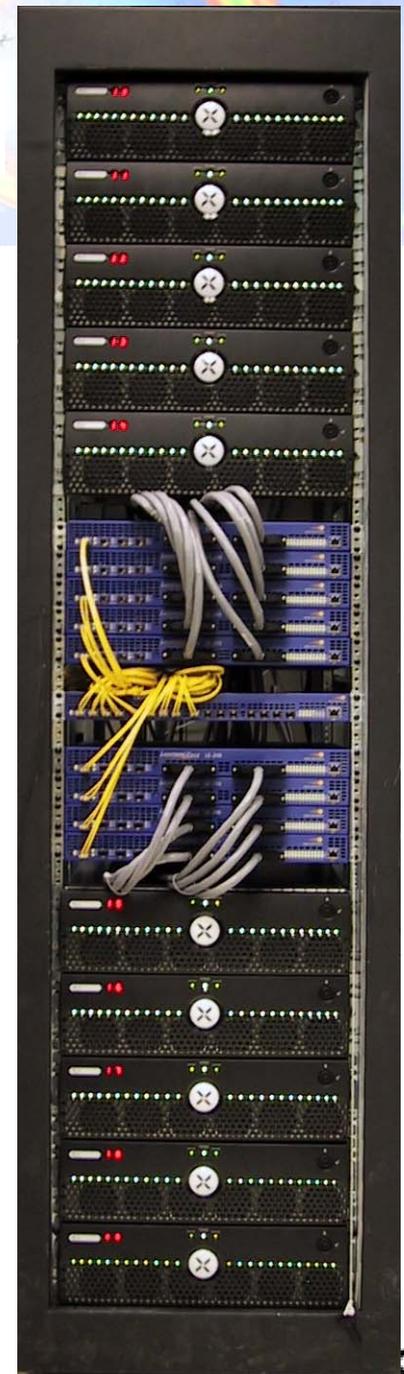
RLX System™ 300ex

- Interchangeable blades
 - Intel, Transmeta, or both.
- Switched-based management

- 3U vertical space
 - 5.25" x 17.25" x 25.2"
- Two hot-pluggable 450W power supplies
 - Load balancing
 - Auto-sensing fault tolerance
- System midplane
 - Integration of system power, management, and network signals.
 - Elimination of internal system cables.
 - Enabling efficient hot-pluggable blades.
- Network cards
 - Hub-based management.
 - Two 24-port interfaces.

"Green Destiny" Bladed Beowulf

- A 240-Node Beowulf in One Cubic Meter
- Each Node
 - ◆ 667-MHz Transmeta TM5600 CPU
 - ☞ Upgrade to 933-MHz Transmeta TM5800 CPUs
 - ◆ 640-MB RAM
 - ◆ 20-GB hard disk
 - ◆ 100-Mb/s Ethernet (up to 3 interfaces)
- Total
 - ◆ 160 Gflops peak (224 Gflops with upgrade)
 - ◆ 240 nodes
 - ◆ 150 GB of RAM (expandable to 276 GB)
 - ◆ 4.8 TB of storage (expandable to 38.4 TB)





Who Cares? So What? It's a Smaller Beowulf ...

- Goal
 - ◆ Improve *efficiency, reliability, and availability* (ERA) in large-scale computing systems.
 - ◆ Reduce the *total cost of ownership* (TCO).

- How to quantify ERA?

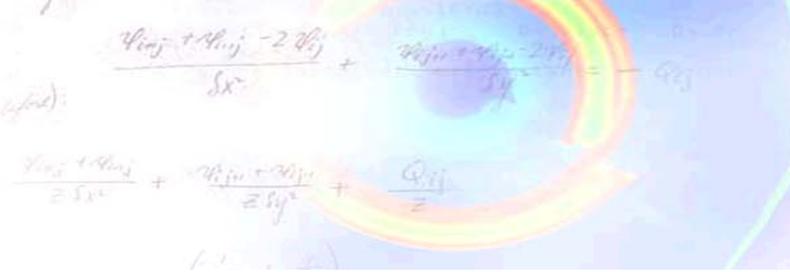
- What exactly is TCO?
 - ◆ Can it be concretely quantified?
 - ◆ Or is it a "foofy" metric?



Performance Metrics



What is TCO?



- Cost of Acquisition ← Fixed, one-time cost
 - ◆ \$\$\$ to buy the supercomputer.
- Cost of Operation ← Variable, recurring cost
 - ◆ Administration
 - ☞ \$\$\$ to build, integrate, configure, maintain, and upgrade the supercomputer over its lifetime.
 - ◆ Power & Cooling
 - ☞ \$\$\$ in electrical power and cooling that is needed to maintain the operation of the supercomputer.
 - ◆ Downtime
 - ☞ \$\$\$ lost due to the downtime (unreliability) of the system.
 - ◆ Space
 - ☞ \$\$\$ spent to house the system.



Total Price-Performance Ratio

- Price-Performance Ratio
 - ◆ *Price = Cost of Acquisition*
 - ◆ *Performance = Floating-Point Operations Per Second*
- Total Price-Performance Ratio (ToPPeR)
 - ◆ *Total Price = Total Cost of Ownership (TCO)*
 - ◆ *Performance = Floating-Point Operations Per Second*



Quantifying TCO?

- Why is TCO hard to quantify?
 - ◆ Components
 - ☞ Acquisition + Administration + Power + Downtime + Space



Quantifying TCO?

- Why is TCO hard to quantify?

- ◆ Components

- ☞ Acquisition + Administration + Power + Downtime + Space

Too Many Hidden Costs
Institution-specific



Quantifying TCO?

- Why is TCO hard to quantify?

- ◆ Components

- ☞ **Acquisition** + Administration + Power + Downtime + Space

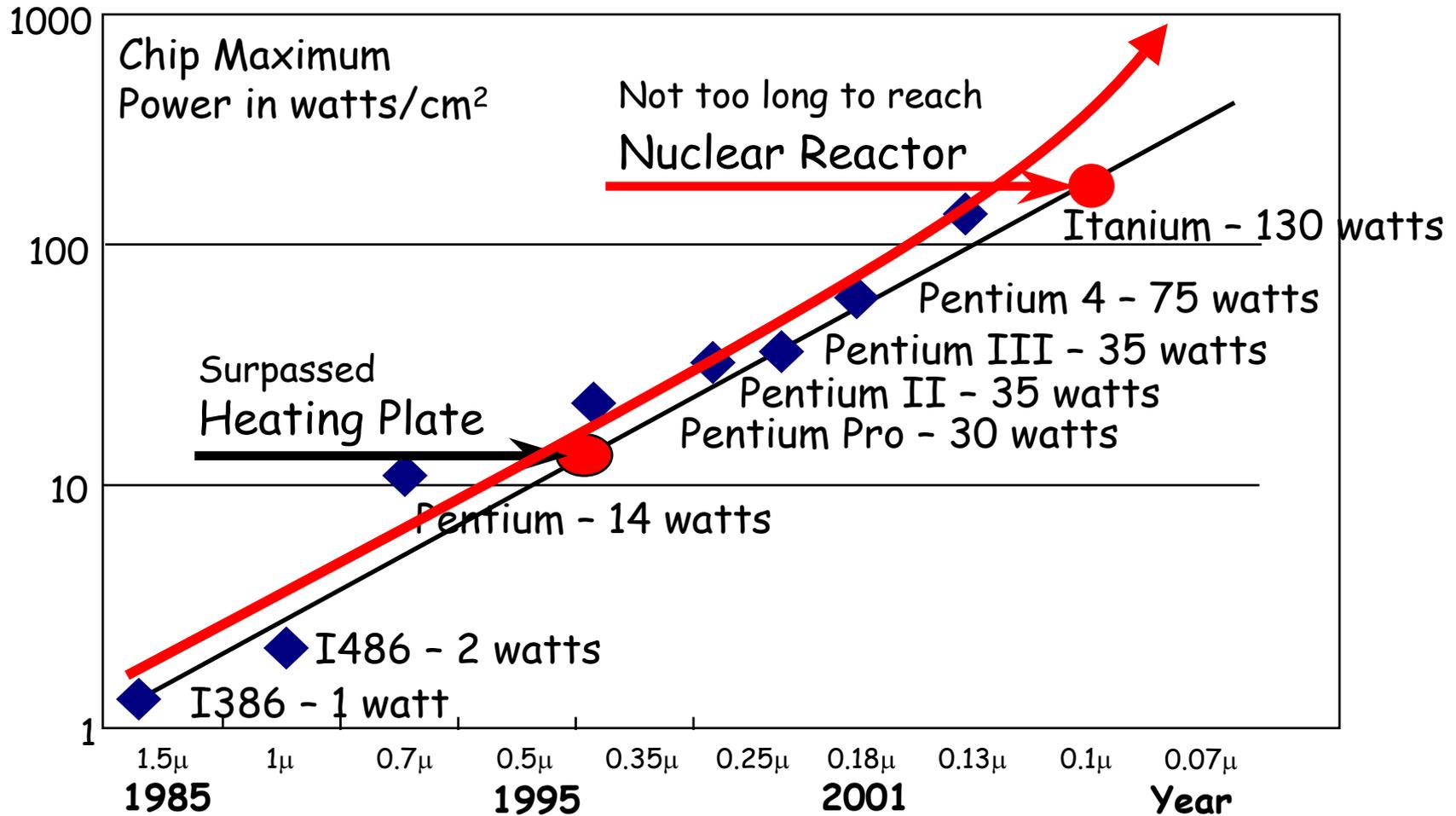
Too Many Hidden Costs
Institution-Specific

- ◆ Traditional Focus: **Acquisition** (i.e., equipment cost)

- ☞ Cost Efficiency: Price/Performance Ratio



Moore's Law for Power Dissipation



Source: Fred Pollack, Intel. New Microprocessor Challenges in the Coming Generations of CMOS Technologies, MICRO32 and Transmeta



Power, Temperature, Reliability

- What's wrong with high power?
 - ◆ Costs \$\$\$ to power such a system; costs \$\$\$ to cool it.
 - ◆ Causes reliability problems. Why?
 - ☞ Higher power implies higher temperatures.
- Arrhenius' Equation (circa 1980s)
 - ◆ As temperature increases by 10°C ...
 - ☞ The failure rate of a system *doubles*.
 - ☞ The reliability of a system is cut in *half*.
 - ◆ Twenty years of unpublished empirical data .



Empirical Data on Temperature

- From off to system boot-up, after 25 seconds:

Processor	Clock Freq.	Voltage	Peak Temp.*
Intel Pentium III-M	500 MHz	1.6 V	252° F (122° C)
Transmeta Crusoe TM5600	600 MHz	1.6 V	147° F (64° C)

*Peak temperature measured with *no* cooling.

- Arrhenius' Equation
 - Every 10° C increase, *doubles* the failure rate.

Implication: Without cooling facilities, PIII-M is 32 times more likely to fail!



Summary of Performance Metrics

- Total Price/Performance Ratio (ToPPeR)
 - ◆ Price is more than the *cost of acquisition*.
 - ◆ Operational costs: sys admin, power & cooling, space, downtime.
- Performance/Power Ratio → "Power Efficiency"
 - ◆ How efficiently does a computing system use energy?
 - ◆ How does this affect reliability and availability?
 - ☞ Higher Power Dissipation α Higher Temperature α Higher Failure Rate
- Performance/Space Ratio → "Space Efficiency"
 - ◆ How efficiently does a computing system use space?
 - ◆ Performance has increased by 2000 since the Cray C90; performance/sq. ft. has only increased by 65.



Benchmark Results



Gravitational Microkernel Benchmark (circa June 2002)

Processor	Math sqrt	Karp sqrt
500-MHz Intel PIII	87.6	137.5
533-MHz Compaq Alpha EV56	76.2	178.5
633-MHz Transmeta TM5600	115.0	144.6
800-MHz Transmeta TM5800	174.1	296.6
375-MHz IBM Power3	298.5	379.1
1200-MHz AMD Athlon MP	350.7	452.5

Units are in Mflops.

Memory Bandwidth for Transmetas (via STREAMS):
350 MB/s

SSS Demo at SC 2001

MetaBlade: 24 ServerBlade 633s

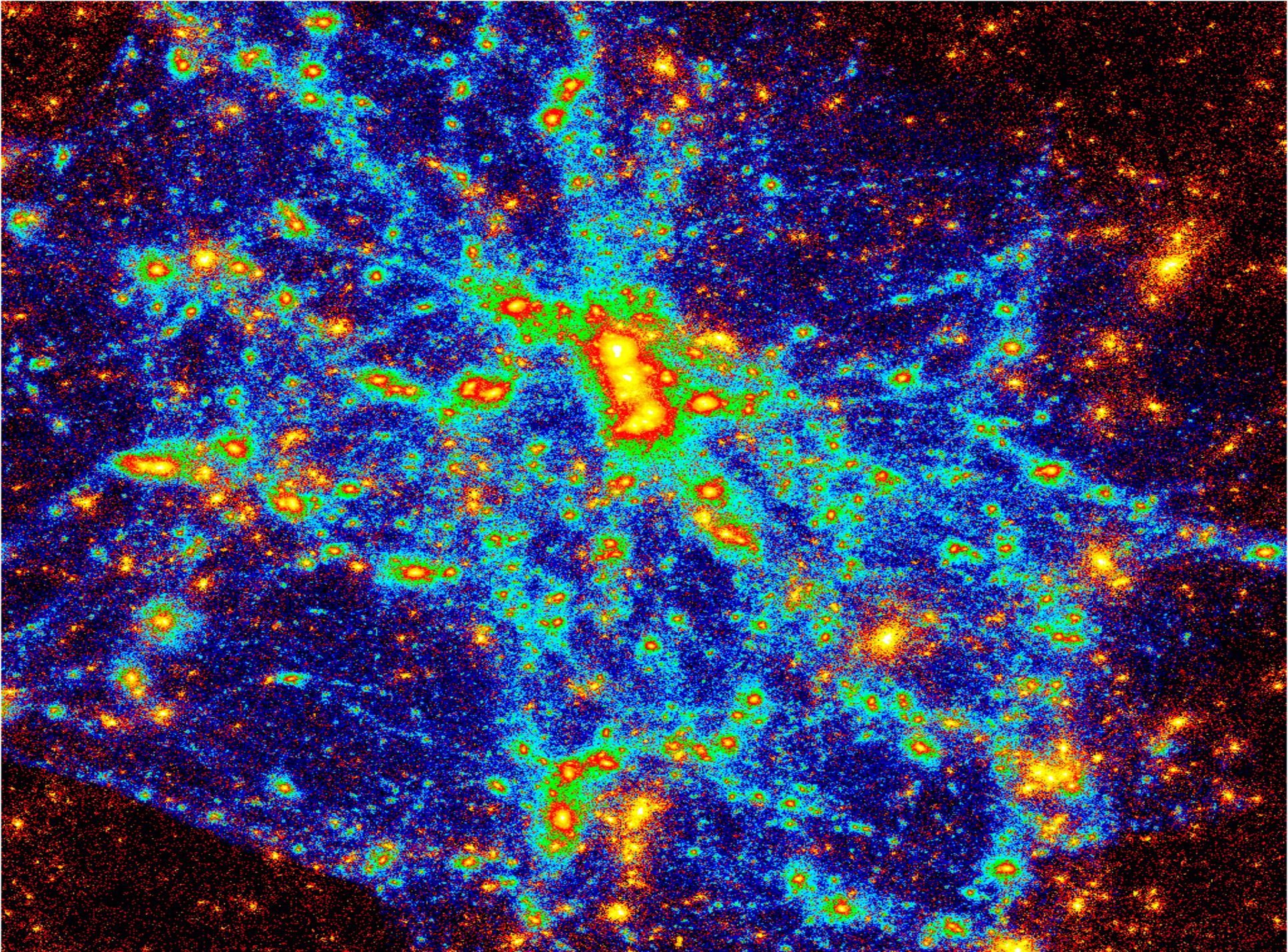


MetaBlade2: 24 ServerBlade 800s



- MetaBlade Bladed Beowulf: 2.1 Gflops (unoptimized)
- MetaBlade2 Bladed Beowulf: 3.3 Gflops (unoptimized)

No failures since September 2001 despite no cooling facilities.





Treecode Benchmark for n-Body

Site	Machine	CPUs	Gflops	Mflops/CPU
NERSC	IBM SP-3	256	57.70	225.0
LANL	SGI O2K	64	13.10	205.0
LANL	Green Destiny	212	38.90	183.5
SC'01	MetaBlade2	24	3.30	138.0
LANL	Avalon	128	16.16	126.0
LANL	Loki	16	1.28	80.0
NASA	IBM SP-2	128	9.52	74.4
SC'96	Loki+Hyglac	32	2.19	68.4
Sandia	ASCI Red	6800	464.90	68.4
CalTech	Naegling	96	5.67	59.1
NRL	TMC CM-5E	256	11.57	45.2



"Cost Efficiency" Metrics

- Price-Performance Ratio
 - ◆ *Price = Cost of Acquisition*
 - ✓ ◆ *Performance = Floating-Point Operations Per Second*
- Total Price-Performance Ratio (ToPPeR)
 - ◆ *Total Price = Total Cost of Ownership (TCO)*
 - ✓ ◆ *Performance = Floating-Point Operations Per Second*



ToPPeR Metric

- ToPPeR: Total Price-Performance Ratio (over the lifetime of a 24-node cluster in a 80° F environment)

Cost Parameter	Alpha	Athlon	PIII	P4	TM5600
Acquisition	\$17K	\$15K	\$16K	\$17K	\$26K
System Admin	\$60K	\$60K	\$60K	\$60K	\$5K
Power & Cooling	\$11K	\$6K	\$6K	\$11K	\$2K
Space	\$8K	\$8K	\$8K	\$8K	\$2K
Downtime	\$12K	\$12K	\$12K	\$12K	\$1K
TCO (four yrs)	\$108K	\$101K	\$102K	\$108K	\$36K

- *Problem: Too many hidden costs & institution-specific*
- ToPPeR metric is approximately 2x better ...



Price/Performance vs. ToPPeR

- Green Destiny
 - ◆ Price/Performance Ratio
 - ☞ \$26K / 38.9 Gflops = \$0.67 / Mflop
 - ◆ ToPPeR (Total Price/Performance Ratio)
 - ☞ \$36K / 38.9 Gflops = \$0.92 / Mflop
- But ToPPeR is a "foofy" metric ...



Parallel Computing Platforms

- Avalon (1996)
 - ◆ 140-Node *Traditional Beowulf Cluster*
- ASCI Red (1996)
 - ◆ 9632-CPU *MPP*
- ASCI White (2000)
 - ◆ 512-Node (8192-CPU) *Cluster of SMPs*
- Green Destiny (2002)
 - ◆ 240-Node *Bladed Beowulf Cluster*



Parallel Computing Platforms Running the N-body Code

Machine	Avalon Beowulf	ASCI Red	ASCI White	Green Destiny
Year	1996	1996	2000	2002
Performance (Gflops)	18	600	2500	39
Area (ft ²)	120	1600	9920	6
Power (kW)	18	1200	2000	5
DRAM (GB)	36	585	6200	150
Disk (TB)	0.4	2.0	160.0	4.8
DRAM density (MB/ft ²)	300	366	625	25000
Disk density (GB/ft ²)	3.3	1.3	16.1	800.0
Power density (watts/ft ²)	150	750	202	833
Space efficiency (Mflops/ft ²)	150	375	252	6500
Power efficiency (Mflops/watt)	1.0	0.5	1.3	7.5



Parallel Computing Platforms Running the N-body Code

Machine	Avalon Beowulf	ASCI Red	ASCI White	Green Destiny
Year	1996	1996	2000	2002
Performance (Gflops)	18	600	2500	39
Area (ft ²)	120	1600	9920	6
Power (kW)	18	1200	2000	5
DRAM (GB)	36	585	6200	150
Disk (TB)	0.4	2.0	160.0	4.8
DRAM density (MB/ft ²)	300	366	625	25000
Disk density (GB/ft ²)	3.3	1.3	16.1	800.0
Power density (watts/ft ²)	150	750	202	833
Space efficiency (Mflops/ft ²)	150	375	252	6500
Power efficiency (Mflops/watt)	1.0	0.5	1.3	7.5



Green Destiny vs. Japanese Earth Simulator

Machine	Green Destiny+	Earth Simulator
Year	2002	2002
LINPACK Performance (Gflops)	144 (ext.)	35,860
Area (ft ²)	6	17,222
Power (kW)	5	7,000
Cost efficiency (\$/Mflop)	2.33	11.15
Space efficiency (Mflops/ft ²)	24,000	2,085
Power efficiency (Mflops/watt)	28.8	5.13

Disclaimer: This is not exactly a fair comparison. Why?

- (1) LINPACK performance is extrapolated for Green Destiny+.
- (2) Use of area and power does *not* scale linearly.

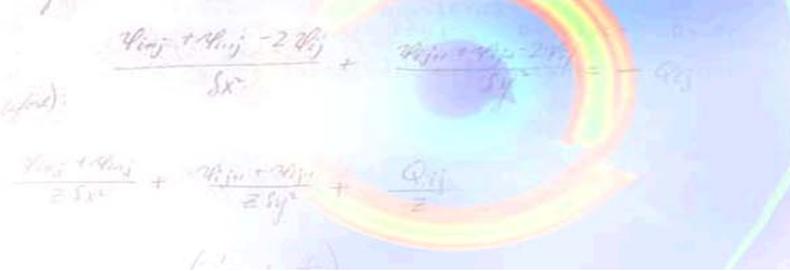


Discussion: Interesting Tidbits

- DARPA Contributes \$2M to IBM's Low Power Center in Aug. 2001
 - ◆ <http://www.computerworld.com/industrytopics/defense/story/0,10801,73289,00.html>
- Transmeta performance on N-body code can match Intel performance on a clock-for-clock-cycle basis.
 - ◆ Problem: Fastest Transmeta? Fastest Intel?
- Low component count on blade server enhances reliability.
 - ◆ 100 parts per RLX node vs. 800-1000 parts per typical node.
- Intel-based Bladed Beowulf: 18 nodes in 3U
 - ◆ 80° F environment: "Silent" failure on LINPACK.
1/3 of nodes inaccessible.
 - ◆ 65 ° F environment: ~20% better performance vs. 933-MHz Transmeta.
- Why 10/100? GigE has been available for two years now.
 - ◆ In 2000-01, GigE ~12-15 W. Now, GigE ~6-8W?
- Systems community vs. applications community.



Status



Recent Work

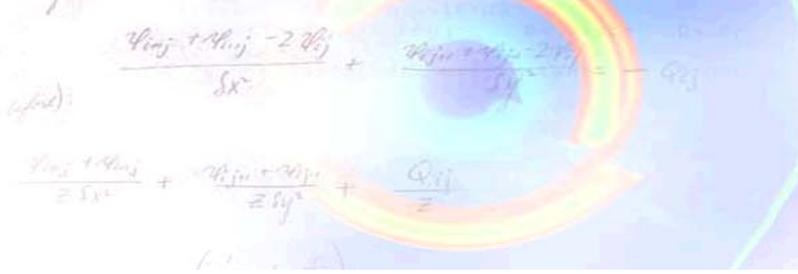
- ◆ April 2002: Assembled and integrated a 240-node Beowulf in one cubic meter called *Green Destiny*.
- ◆ July 2002: Worked with Transmeta to demonstrate comparable performance to similarly-clocked Intels.
- ◆ July 2002: Worked with DOE SciDAC-funded *3-D Supernova* project to demo "base code" on *Green Destiny*. (A vertically-integrated solution from hardware on up to the application.)
- ◆ August 2002: Completed mpiBLAST code. Presented at IEEE Bioinformatics. Demonstrated super-linear speed-up.

Future Work

- ◆ Demo first 3-D supernova on a Linux-based cluster at SC.
- ◆ Work with additional code teams, e.g., climate modeling, computational fluid dynamics, large-scale molecular dynamics.
- ◆ Upgrade *Green Destiny* processors from 667 MHz to 933 MHz.



Conclusion



- New Performance Metrics
 - ◆ Overall Efficiency
 - ☞ ToPPeR: Total Price-Performance Ratio
 - ◆ Power Efficiency
 - ☞ Performance-Power Ratio
 - ◆ Space Efficiency
 - ☞ Performance-Space Ratio
- Predictions
 - ◆ Traditional clustering and supercomputing as we know it will *NOT* scale to petaflop computing due to issues of efficiency, reliability, and availability.
 - ◆ "Supercomputing in Small Spaces" is a single step in the right direction ...



Conclusion

- Keeping It In Perspective
 - ◆ The “Supercomputing in Small Spaces” project (<http://sss.lanl.gov>) is *not* meant to replace today’s large supercomputers.
 - ☞ Focus on metrics related to efficiency, reliability, and availability (ERA) rather than raw performance.
 - i.e., SSS = “Toyota Camry” of supercomputing.
 - ☞ Works particularly well as a departmental cluster (or even institutional cluster if there exists power and space constraints).



Acknowledgments

- Technical Co-Leads
 - ◆ Mike Warren and Eric Weigle
- Contributions
 - ◆ Mark Gardner, Adam Engelhart, Gus Hurwitz
- Enablers
 - ◆ J. Thorp, A. White, R. Oldehoeft, and D. Lora (LACSI)
 - ◆ W. Feiereisen and S. Lee (CCS Division Office)
- Funding Agencies
 - ◆ LACSI
 - ◆ IA-Linux
- Encouragement & Support
 - ◆ Gordon Bell, Chris Hipp, Linus Torvalds



The "Hype": A Sampling of Press Coverage

- "Not Your Average Supercomputer," *Communications of the ACM*, 8/02.
- "At Los Alamos, Two Visions of Supercomputing," *The New York Times*, 6/25/02.
- "Two Directions for the Future of Supercomputing," *slashdot.org*, 6/25/02.
- "Researchers Deliver Supercomputing in Smaller Package," *Yahoo! Finance*, 6/4/02.
- "Computer World Faces Heat Wave," *Santa Fe New Mexican*, 6/3/02.
- "Supercomputing Coming to a Closet Near You?" *HPCwire*, 5/31/02.
- "Smaller, Slower Supercomputers May Someday May Win The Race," *HPCwire*, 5/31/02.
- "Supercomputing Coming to a Closet Near You?" *PCworld.com*, 5/27/02.
- "Bell, Torvalds Usher Next Wave of Supercomputing," *HPCwire*, 5/24/02.
- "Supercomputing Cut Down to Size," *Personal Computer World*, 5/22/02.
- "Bell, Torvalds Usher Next Wave of Supercomputing," *CNN.com*, 5/21/02.
- "Transmeta's Low Power Finds Place in Supercomputers," *ZDNet News*, 5/20/02.
- "Transmeta Blades Power Landmark Supercomputer Breakthrough," *The Register*, 5/20/02.



SUPERCOMPUTING
in **SMALL SPACES**

<http://sss.lanl.gov>

Wu-chun (Wu) Feng

Research and Development in Advanced Network Technology



<http://www.lanl.gov/radiant>